

# Stick-Breaking Policy Learning in Dec-POMDPs

Miao Liu<sup>1</sup>, Christopher Amato<sup>2</sup>, Xuejun Liao<sup>3</sup>, Lawrence Carin<sup>3</sup>, and Jonathan P. How<sup>4</sup>

<sup>1</sup>Laboratory for Information and Decision Systems , MIT, Cambridge, MA 02139, USA, miaoliu@mit.edu

<sup>2</sup>Department of Computer Science , University of New Hampshire, Durham, NH 03824, USA , camato@cs.unh.edu

<sup>3</sup>Department of Electrical & Computer Engineering , Duke University, Durham NC 27708, USA , {xjliao,lcarin}@duke.edu

<sup>4</sup>Department of Aeronautics & Astronautics , MIT, Cambridge, MA 02139, USA , jhow@mit.edu

## Abstract

Expectation maximization (EM) has recently been shown to be an efficient algorithm for learning finite-state controllers (FSCs) in large decentralized POMDPs (Dec-POMDPs). However, current methods use fixed-size FSCs and often converge to maxima that are far from optimal. This paper considers a variable-size FSC to represent the local policy of each agent. These variable-size FSCs are constructed using a stick-breaking prior, leading to a new framework called *decentralized stick-breaking policy representation* (Dec-SBPR). This approach learns the controller parameters with a variational Bayesian algorithm without having to assume that the Dec-POMDP model is available. The performance of Dec-SBPR is demonstrated on several benchmark problems, showing that the algorithm scales to large problems while outperforming other state-of-the-art methods.

## 1 Introduction

Decentralized partially observable Markov decision processes (Dec-POMDPs) [3, 25] provide a general framework for solving the cooperative multiagent sequential decision-making problems that arise in numerous applications, including robotic soccer [24], transportation [4], extraplanetary exploration [8], and traffic control [35]. Dec-POMDPs can be viewed as a POMDP controlled by multiple distributed agents. These agents make decisions based on their own local streams of information (i.e., observations), and their joint actions control the global state dynamics and the expected reward of the team. Because of the decentralized decision-making, an individual agent generally does not have enough information to compute the global belief state, which is a

sufficient statistic for decision making in POMDPs. This makes generating an optimal solution in a Dec-POMDP more difficult than for a POMDP [10], especially for long planning horizons.

To circumvent the difficulty of solving long-horizon Dec-POMDPs optimally, while still generating a high quality policy, this paper presents scalable learning methods using a finite memory policy representation. For infinite-horizon problems (which continue for an infinite number of steps), significant progress has been made with agent policies represented as finite-state controllers (FSCs) that map observation histories to actions [9, 2]. Recent work has shown that expectation-maximization (EM) [14] is a scalable method for generating controllers for large Dec-POMDPs [19, 27]. In addition, EM has also been shown to be an efficient algorithm for policy-based reinforcement learning (RL) in Dec-POMDPs, where agents learn FSCs based on trajectories, without knowing or learning the Dec-POMDP model [35].

An important and yet unanswered question is how to define an appropriate number of nodes in each FSC. Previous work assumes a fixed FSC size for each agent, but the number of nodes affects both the quality of the policies and the convergence rate. When the number of nodes is too small, the FSC is unable to represent the optimal policy and therefore will quickly converge to a sub-optimal result. By contrast, when the number is too large, the FSC overfits data, often yielding slow convergence and, again, a sub-optimal policy.

This paper uses a Bayesian nonparametric approach to determine the appropriate controller size in a variable-size FSC. Following previous methods [35, 25], learning is assumed to be centralized, and execution is decentralized. That is, learning is accomplished offline based on all available information, but the optimization is only over decentralized solutions. Such a controller is constructed using the stick-breaking (SB) prior [18]. The SB prior allows the number of nodes to be variable, but the set of nodes that is actively used by the controller is encouraged to be compact. The nodes that are actually used are determined by the posterior, combining the SB prior and the information from trajectory data. The framework is called the *decentralized stick-breaking policy representation* (Dec-SBPR) to recognize the role of the SB prior.

In addition to the use of variable-size FSCs, the paper also makes several other contributions. Specifically, our algorithm directly operates on the (shifted) empirical value function of Dec-POMDPs, which is simpler than the likelihood functions (a mixture of dynamic Bayes nets (DBNs)) in existing planning-as-inference frameworks [19, 35]. Moreover, we derive a variational Bayesian (VB) algorithm for learning the Dec-SBPR based only on the agents’ trajectories (or episodes) of actions, observations, and rewards. The VB algorithm is linear in the number of agents and at most square in the problem size, and is therefore scalable to large application domains. In practice, these trajectories can be generated by a simulator or a set of real-world experiences that are provided, and this batch data scenario is general and realistic, as it is widely adopted in learning from demonstration [23], and reinforcement learning. To the best of our knowledge, this is the first application of Bayesian nonparametric methods to the difficult and little-studied problem of policy-based RL in Dec-POMDPs, and the proposed method is able to generate high-quality solutions for large problems.

## 2 Background and Related Work

Before introducing the proposed method, we first describe the Dec-POMDP model and some related work.

## 2.1 Decentralized POMDPs

A Dec-POMDP can be represented as a tuple  $\mathcal{M} = \langle \mathcal{N}, \mathcal{A}, \mathcal{S}, \mathcal{O}, \mathcal{T}, \Omega, \mathcal{R}, \gamma \rangle$ , where  $\mathcal{N} = \{1, \dots, N\}$  is a finite set of agent indices;  $\mathcal{A} = \otimes_n \mathcal{A}_n$  and  $\mathcal{O} = \otimes_n \mathcal{O}_n$  respectively are sets of joint actions and observations, with  $\mathcal{A}_n$  and  $\mathcal{O}_n$  available to agent  $n$ . At each step, a joint action  $\vec{a} = (a_1, \dots, a_N) \in \mathcal{A}$  is selected and a joint observation  $\vec{o} = (o_1, \dots, o_N)$  is received;  $\mathcal{S}$  is a set of finite world states;  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the state transition function with  $\mathcal{T}(s'|s, \vec{a})$  denoting the probability of transitioning to  $s'$  after taking joint action  $\vec{a}$  in  $s$ ;  $\Omega : \mathcal{S} \times \mathcal{A} \times \mathcal{O} \rightarrow [0, 1]$  is the observation function with  $\Omega(\vec{o}|s', \vec{a})$  the probability of observing  $\vec{o}$  after taking joint action  $\vec{a}$  and arriving in state  $s'$ ;  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function with  $r(s, \vec{a})$  the immediate reward received after taking joint action  $\vec{a}$  in  $s$ ;  $\gamma \in [0, 1)$  is a discount factor. A global reward signal is generated for the team of agents after joint actions are taken, but each agent only observes its local observation. Because each agent lacks access to other agents' observations, each agent maintains a local policy  $\pi_n$ , defined as a mapping from local observation histories to actions. A joint policy consists of the local policies of all agents. For an infinite-horizon Dec-POMDP with initial belief state  $b_0$ , the objective is to find a joint policy  $\Psi = \otimes_n \Psi_n$ , such that the value of  $\Psi$  starting from  $b_0$ ,  $V^\Psi(b(s_0)) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, \vec{a}_t) | b_0, \Psi]$ , is maximized.

An FSC is a compact way to represent a policy as a mapping from histories to actions. Formally, a stochastic FSC for agent  $n$  is defined as a tuple  $\Theta_n = \langle \mathcal{A}_n, \mathcal{O}_n, \mathcal{Z}_n, \mu_n, W_n, \pi_n \rangle$ , where,  $\mathcal{A}_n$  and  $\mathcal{O}_n$  are the same as defined in the Dec-POMDP;  $\mathcal{Z}_n$  is a finite set of controller nodes for agent  $n$ ;  $\mu_n$  is the initial node distribution with  $\mu_n^z$  the probability of agent  $n$  initially being in  $z$ ;  $W_n$  is a set of Markov transition matrices with  $W_{n,a,o}^{z,z'}$  denoting the probability of the controller node transiting from  $z$  to  $z'$  when agent  $n$  takes action  $a$  in  $z$  and sees observation  $o$ ;  $\pi_n$  is a set of stochastic policies with  $\pi_{n,z}^a$  the probability of agent  $n$  taking action  $a$  in  $z$ .

For simplicity, we use the following notational conventions.  $\mathcal{Z}_n = \{1, 2, \dots, C_n\}$ , where  $C_n \stackrel{\text{def.}}{=} |\mathcal{Z}_n|$  is the cardinality of  $\mathcal{Z}_n$ , and  $\mathcal{A}_n$  and  $\mathcal{O}_n$  follow similarly.  $\Theta = \{\Theta_1, \dots, \Theta_N\}$  is the joint FSC of all agents. A consecutively-indexed variable is abbreviated as the variable with the index range shown in the subscript or superscript; when the index range is obvious from the context, a simple “.” is used instead. Thus,  $a_{n,0:T} = (a_{n,0}, a_{n,1}, \dots, a_{n,T})$  represents the actions of agent  $n$  from step 0 to  $T$  and  $W_{n,a,o}^{z,:} = (W_{n,a,o}^{z,1}, W_{n,a,o}^{z,2}, \dots, W_{n,a,o}^{z,|\mathcal{Z}_n|})$  represents the node transition probabilities for agent  $n$  when starting in node  $z$ , taking action  $a$  and seeing observation  $o$ . Given  $h_{n,t} = \{a_{n,0:t-1}, o_{n,1:t}\}$ , a local history of actions and observations up to step  $t$ , as well as an agent controller,  $\Theta_n$ , we can calculate a local policy  $p(a_{n,t} | h_{n,t}, \Theta_n)$ , the probability that agent  $n$  chooses its action  $a_{n,t}$ .

## 2.2 Planning as Inference in Dec-POMDPs

A Dec-POMDP planning problem can be transformed into an inference problem and then efficiently solved by EM algorithms. The validity of this method is based on the fact that by introducing binary rewards  $R$  such that  $P(R = 1 | a, s) \propto r(a, s)$ ,  $\forall a \in \mathcal{A}, s \in \mathcal{S}$  and choosing the geometric time prior  $p(T) = \gamma^T(1-\gamma)$ , maximizing the likelihood  $L(\Theta) = P(R = 1; \Theta) = \sum_{T=0}^{\infty} P(T)P(R = 1 | T; \Theta)$  of a mixture of dynamic Bayes nets is equivalent to optimizing the associated Dec-POMDP policy,

as the joint-policy value  $V(\Theta)$  and  $L(\Theta)$  can be related through an affine transform [19]

$$L(\Theta) = \frac{(1 - \gamma)(V(\Theta) - \sum_T \gamma^T R_{min})}{R_{max} - R_{min}} = \frac{(1 - \gamma)\hat{V}(\Theta)}{R_{max} - R_{min}}, \quad (1)$$

where  $R_{max}$  and  $R_{min}$  are maximum and minimum reward values and  $\hat{V}(\Theta) \stackrel{def}{=} V(\Theta) - \sum_T \gamma^T R_{min}$  is a shifted value.

Previous EM methods [19, 27] have achieved success in scaling to larger problems by factoring the distribution over states and histories for inference, but these methods require using a Dec-POMDP model to construct a Bayes net for policy evaluation. When the exact model parameters  $\mathcal{T}$ ,  $\Omega$  and  $\mathcal{R}$  are unknown, one needs to solve a reinforcement learning (RL) problem. To address this important yet less addressed problem, a global empirical value function extended from the single-agent case [20], is constructed based on all the action, observation and reward trajectories, and the product of local policies of all agents. This serves as the basis for learning (fixed-size) FSCs in RL settings.

**Definition 1.** (*Global empirical value function*) Let  $\mathcal{D}^{(K)} = \{(\vec{o}_0^k \vec{a}_0^k r_0^k \vec{o}_1^k \vec{a}_1^k r_1^k \dots \vec{o}_{T_k}^k \vec{a}_{T_k}^k r_{T_k}^k)\}_{k=1, \dots, K}$  be a set of episodes resulting from  $N$  agents who choose actions according to  $\Psi = \otimes_n \Psi_n$ , a set of stochastic behavior policies with  $p^{\Psi_n}(a|h) > 0, \forall$  action  $a, \forall$  history  $h$ . The global empirical value function is defined as  $\hat{V}(\mathcal{D}^{(K)}; \Theta) \stackrel{def}{=} \sum_{k=1}^K \sum_{t=0}^{T_k} \frac{\gamma^t (r_t^k - R_{min})}{K} \frac{\prod_{\tau=0}^t \prod_{n=1}^N p(a_{n,\tau}^k | h_{n,\tau}^k, \Theta_n)}{\prod_{\tau=0}^t \prod_{n=1}^N p^{\Psi_n}(a_{n,\tau}^k | h_{n,\tau}^k)}$  where  $h_{n,t}^k = (a_{n,0:t-1}^k, o_{n,1:t}^k)$ ,  $0 \leq \gamma < 1$  is the discount.

According to the strong law of large numbers [32],  $\hat{V}(\Theta) = \lim_{K \rightarrow \infty} \hat{V}(\mathcal{D}^{(K)}; \Theta)$ , i.e., with a large number of trajectories, the empirical value function  $\hat{V}(\mathcal{D}^{(K)}; \Theta)$  approximates  $\hat{V}(\Theta)$  accurately. Hence, applying (1),  $\hat{V}(\mathcal{D}^{(K)}; \Theta)$  approximates  $L(\Theta)$ , and offers an objective for learning the decentralized policies and can be directly maximized by the EM algorithms in [20].

### 3 Bayesian Learning of Policies

EM algorithms infer policies based on fixed-size representation and observed data only, it is difficult to explicitly handle model uncertainty and encode prior (or expert) knowledge. To address these issues, a Bayesian learning method is proposed in this section. This is accomplished by measuring the likelihood of  $\Theta$  using  $L(\mathcal{D}^{(K)}; \Theta)$ , which is combined with the prior  $p(\Theta)$  in Bayes' rule to yield the posterior

$$p(\Theta | \mathcal{D}^{(K)}) = L(\mathcal{D}^{(K)}; \Theta) p(\Theta) [p(\mathcal{D}^{(K)})]^{-1}, \quad (2)$$

where  $p(\mathcal{D}^{(K)})$  is the marginal likelihood of the joint FSC and, up to additive constant, proportional to the marginal value function,

$$\begin{aligned} \hat{V}(\mathcal{D}^{(K)}) &\stackrel{def}{=} \int \hat{V}(\mathcal{D}^{(K)}; \Theta) p(\Theta) d\Theta \\ &\propto \int L(\mathcal{D}^{(K)}; \Theta) p(\Theta) d\Theta = p(\mathcal{D}^{(K)}). \end{aligned} \quad (3)$$

To compute the posterior,  $p(\Theta | \mathcal{D}^{(K)})$ , Markov chain Monte Carlo (MCMC) simulation [32] is the most straight forward method. However, MCMC is costly in terms of computation and storage, and

lacks a strong convergence guarantee. An alternative is a variational Bayes (VB) method [7], which performs approximate posterior inference by minimizing the Kullback-Leibler (KL) divergence between the true and approximate posterior distributions. Because the VB method has a (local) convergence guarantee and is able to trade-off scalability and accuracy, we focus on the derivation of VB method here. Denoting  $q(\Theta)$  as the variational approximation to  $p(\Theta|\mathcal{D}^{(K)})$ , and  $q_t^k(\vec{z}_{0:t}^k)$  as the approximation to  $p(\vec{z}_{0:t}^k|\vec{o}_{1:t}^k, \Theta)$ , a VB objective function <sup>1</sup> is

$$\text{KL}(\{q_t^k(\vec{z}_{0:t}^k)q(\Theta)\}_{k=1:K}||\{\nu_t^k p(\vec{z}_{0:t}^k, \Theta)\}_{k=1:K}) = \ln \hat{V}(\mathcal{D}^{(K)}) - \text{LB}(\{q_t^k(\vec{z}_{0:t}^k)\}, q(\Theta)), \quad (4)$$

where

$$\text{LB}(\{q_t^k(\vec{z}_{0:t}^k)\}, q(\Theta)) \stackrel{\text{def.}}{=} \sum_{k,t,z_{1:N,0:t}^k} \int \frac{q_t^k(\vec{z}_{0:t}^k)}{K/q(\Theta)} \ln \frac{\nu_t^k p(\vec{z}_{0:t}^k, \Theta|h_{0:t}^k)}{q_t^k(\vec{z}_{0:t}^k)q(\Theta)} d\Theta \quad (5)$$

is the lower bound of  $\ln \hat{V}(\mathcal{D}^{(K)})$  and

$$\nu_t^k \stackrel{\text{def.}}{=} \frac{\gamma^t r_t^k \prod_{n=1}^N p(a_{n,0:t}^k | o_{n,1:t}^k)}{\prod_{n=1}^N \prod_{\tau=0}^t p^{\Psi_n}(a_{n,\tau}^k | h_{n,\tau}^k) \hat{V}(\mathcal{D}^{(K)})}, \forall t, k, \quad (6)$$

is the re-weighted reward. Since  $\ln \hat{V}(\mathcal{D}^{(K)})$  in equation (4) is independent of  $\Theta$  and  $\{q_t^k(\vec{z}_{0:t}^k)\}$ , minimizing the KL divergence is equivalent to maximizing the lower bound, leading to the following constrained optimization problem,

$$\begin{aligned} & \max_{\{q_t^k(\vec{z}_{0:t}^k)\}_{q(\Theta)}} \text{LB}(\{q_t^k(\vec{z}_{0:t}^k)\}, q(\Theta)) \\ & \text{subject to: } q_t^k(\vec{z}_{0:t}^k, \Theta) = \prod_{n=1}^N q_t^k(z_{n,0:t}^k)q(\Theta_n), \\ & \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{z_{1:N,0:t}^k} q_t^k(\vec{z}_{0:t}^k) = K, \quad q_t^k(\vec{z}_{0:t}^k) \geq 0, \forall \vec{z}_{0:t}^k, t, k, \\ & \int p(\Theta) d\Theta = 1 \quad \text{and} \quad p(\Theta) \geq 0, \forall \Theta, \end{aligned} \quad (7)$$

where the constraint in the second line arises both from the mean-field approximation and from the decentralized policy representation, and the last two lines summarize the normalization constraints. It is worth emphasizing that we developed this variational mean-field approximation to optimize a decentralized policy representation, showing that the VB learning problem formulation (7) is both a general and accurate method for the multiagent problem considered in this paper.

### 3.1 Stick-breaking Policy Priors

To solve the Bayesian learning problem described above and obtain the variable-size FSCs, the stick-breaking prior is used to specify the policy's structure. As such, Dec-SBPR is formally given in definition 2.

<sup>1</sup>Refer to the appendix for derivation details

**Definition 2.** The decentralized stick breaking policy representation (Dec-SBPR) is a tuple  $(\mathcal{N}, \mathcal{A}, \mathcal{O}, \mathcal{Z}, \mu, \eta, \rho)$ , where  $\mathcal{N}, \mathcal{A}$  and  $\mathcal{O}$  are as in the definition of Dec-POMDP;  $\mathcal{Z}$  is an unbounded set of nodes indexed by positive integers; for notational simplicity<sup>2</sup>,  $\mu$  are assumed to be deterministic with  $\mu_n^1 = 1, \mu_n^{2:\infty} = 0, \forall n$ ;  $(\eta, \rho)$  determine  $(W, \pi)$ , the FSC parameters defined in section 2.1, as follows

$$W_{n,a,o}^{i,1:\infty} \sim \text{SB}(\sigma_{n,a,o}^{i,1:\infty}, \eta_{n,a,o}^{i,1:\infty}), \quad \pi_{n,i}^{1:|\mathcal{A}_n|} \sim \text{Dir}(\rho_{n,i}^{1:|\mathcal{A}_n|}) \quad (8)$$

where  $\text{Dir}$  represents Dirichlet distribution and  $\text{SB}$  represents the stick-breaking process with  $W_{n,a,o}^{i,j} = V_{n,a,o}^{i,j} \prod_{m=1}^{j-1} (1 - V_{n,a,o}^{i,m})$  and  $V_{n,a,o}^{i,j} \sim \text{Beta}(\sigma_{n,a,o}^{i,j}, \eta_{n,a,o}^{i,j})$ ,  $\eta_{n,a,o}^{i,j} \sim \text{Gamma}(c, d)$ ,  $n = 1, \dots, N$  and  $i, j = 1, \dots, \infty$ .

DECSBPR differs from previous nonparametric Bayesian RL methods [21, 16]. Specifically, Dec-SBPR performs policy-based RL and generalizes the nonparametric Bayesian policy representation of POMDPs [21] to the decentralized domain. Whereas [16] is a model-based RL method that doesn't assume knowledge about the world's model, but explicitly learns it and then performs planning. Moreover, Dec-SBPR further distinguishes from previous methods [16, 21] by the prior distributions and inference methods employed. These previous methods employed hierarchical Dirichlet processes hidden Markov models (HDP-HMM) to infer the number of controller nodes. However, due to the lack of conjugacy between two levels of DPs in the HDP-HMM, a *fully conjugate Bayesian* variational inference does not exist<sup>3</sup>. Therefore, these methods used MCMC which requires high computational and storage costs, making them not ideal for solving large problems. In contrast, Dec-SBPR employs single layer SB priors over FSC transition matrices  $W$  and sparse Gamma priors over SB weight hyperparameters  $\eta$  to bias transition among nodes with smaller indices. A similar framework has been explored to infer HMMs, and we refer readers to [26] for more details.

It is worth noting that SB processes subsume Dirichlet Processes (DPs) [17] as a special case, when  $\sigma_{n,a,o}^{i,j} = 1, \forall i, j, n, a, o$  (in Dec-SBPR). The purpose of using SB priors is to encourage a small number of FSC nodes. Compared to a DP, the SB priors can represent richer patterns of sparse transition between the nodes of an FSC, because it allows arbitrary correlation between the stick-breaking weights (the weights are always negatively correlated in a DP).

### 3.2 Variational Stick-breaking Policy Inference

It is shown in [18] that the random weights constructed by the SB prior are equivalently governed by a generalized Dirichlet distribution (GDD) and are therefore conjugate to the multinomial distribution; hence an efficient variational Bayesian algorithm for learning the decentralized policies can be derived. To accommodate an unbounded number of nodes, we apply the retrospective representation of SB priors [28] to the Dec-SBPR. For agent  $n$ , the SB prior is set with a truncation level  $|\mathcal{Z}_n|$ , taking into account the current occupancy as well as additional nodes reserved for future new occupancies. The solution to (7) under the stick-breaking priors is given in Theorem 3, the proof of which is provided in the appendix.

<sup>2</sup>Nonparametric priors over  $\mu$  can also be used.

<sup>3</sup>The VB method in [12] imposes point-mass proposals over top level DPs, lacking a uncertainty measure.

---

**Algorithm 1: Batch VB Inference for Dec-SBPR**


---

```

1: Input: Episodes  $\mathcal{D}^{(K)}$ , the number of agents  $N$ , initial policies  $\Theta$ , VB lower bound  $\text{LB}_0 = -\text{Inf}$ ,  $\Delta\text{LB} = 1$ ,
   Iter = 0;
2: while  $\Delta\text{LB} > 10^{-3}$  do
3:   for  $k = 1$  to  $K$ ,  $n = 1$  to  $N$  do
4:     Update the global rewards  $\{\hat{\nu}_t^k\}$  using (9).
5:     Compute  $\{\alpha_{\tau}^{n,k}\}$  and  $\{\beta_{t,\tau}^{n,k}\}$ .
6:   end for
7:   Iter = Iter + 1.
8:   Compute  $\text{LB}_{\text{Iter}}$  using (5)
9:    $\Delta\text{LB} = (\text{LB}_{\text{Iter}} - \text{LB}_{\text{Iter}-1})/|\text{LB}_{\text{Iter}-1}|$ 
10:  for  $n = 1$  to  $N$  do
11:    Compute  $\{\xi_{t,\tau}^{n,k}(i, j)\}$  and  $\{\phi_{t,\tau}^{n,k}(i)\}$  using (11).
12:    Update the hyper-parameters of  $\Theta_n$  using (10).
13:    Compute  $|\mathcal{Z}_n|$  using (13).
14:  end for
15: end while
16: Return: Policies  $\{\Theta_n\}_{n=1}^N$ , and controller sizes  $\{|\mathcal{Z}_n|\}_{n=1}^N$ .

```

---

**Theorem 3.** Let  $p(\Theta)$  be constructed by the SB priors defined in (8) with hyper-parameters  $(\hat{\sigma}, \hat{\eta}, \hat{\rho})$ , then iterative application of the following updates leads to monotonic increase of (5), until convergence to a maxima. The updates of  $\{q_t^k\}$  are

$$q_t^k(z_{n,0:t}^k) = \hat{\nu}_t^k p(z_{n,0:t}^k | o_{n,1:t}^k, a_{n,0:t}^k, \tilde{\Theta}_n), \forall n, t, k, \quad (9)$$

where  $\hat{\nu}_t^k$  is computed using (6) with  $\Theta$  replaced by  $\tilde{\Theta} = \{\tilde{\pi}, \tilde{\mu}, \tilde{W}\}$ , a set of under-normalized probability (mass) functions, with  $\tilde{\pi}_{n,i}^a = e^{\langle \ln \pi_{n,i}^a \rangle_{p(\pi|\hat{\rho})}}$ , and  $\tilde{W}_{n,a,o}^{i,j} = e^{\langle \ln W_{n,a,o}^{i,j} \rangle_{p(W|\hat{\sigma}, \hat{\eta})}}$ , and  $\langle \cdot \rangle_p$  denotes expectations of  $\cdot$  with respect to distributions  $p$ . The hyper-parameters of the posterior distribution are updated as

$$\begin{aligned} \hat{\sigma}_{n,a,o}^{i,j} &= \sigma_{n,a,o}^{i,j} + \zeta_{n,a,o}^{i,j}, \quad \hat{\eta}_{n,a,o}^{i,j} = \eta_{n,a,o}^{i,j} + \sum_{l=j+1}^{|\mathcal{Z}_n|} \zeta_{n,a,o}^{i,l}, \\ \hat{\rho}_{n,i}^a &= \rho_{n,i}^a + \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\tau=1}^t \frac{\hat{\nu}_t^k}{K} \phi_{t,\tau-1}^{n,k}(i) \mathbb{I}_a(a_{\tau}^k) \end{aligned} \quad (10)$$

with  $\zeta_{n,a,o}^{i,j} = \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\tau=1}^t \frac{\hat{\nu}_t^k}{K} \xi_{t,\tau-1}^{n,k}(i, j) \mathbb{I}_{a,o}(a_{\tau-1}^k, o_{\tau}^k)$ , where  $\mathbb{I}(\cdot)$  is the indicator function, and both  $\xi_{t,\tau}^{n,k}$  and  $\phi_{t,\tau}^{n,k}$  are marginals of  $q_t^k(z_{n,0:t}^k)$ , i.e.

$$\xi_{t,\tau}^{n,k}(i, j) = p(z_{n,\tau}^k = i, z_{n,\tau+1}^k = j | a_{n,0:t}^k, o_{n,1:t}^k, \tilde{\Theta}_n) \quad (11)$$

$$\phi_{t,\tau}^{n,k}(i) = p(z_{n,\tau}^k = i | a_{n,0:t}^k, o_{n,1:t}^k, \tilde{\Theta}_n) \quad (12)$$

The update equations in Theorem 3 constitute the VB algorithm for learning a variable-size joint FSCs under SB priors with batch data. In particular, (9) is a policy-evaluation step where the rewards are reweighted to reflect the improved marginal value of the new policy posterior updated in the previous iteration, and (10) is a policy-improvement step where the reweighted rewards are

used to further improve the policy posterior. Both steps require (11), which are computed based on  $\alpha_{\tau}^{n,k}(i) = p(z_{n,\tau}^k = i | a_{n,0:\tau}^k, o_{n,1:\tau}^k, \tilde{\Theta}_n)$  and  $\beta_{t,\tau}^{n,k}(i) = \frac{p(a_{n,\tau+1:t}^k | z_{n,\tau}^k = i, o_{n,\tau+1:t}^k, \tilde{\Theta}_n)}{\prod_{\tau'=t}^{\tau} p(a_{n,\tau'}^k | h_{n,\tau'}^k, \tilde{\Theta}_n)}$ ,  $\forall n, k, t, \tau$ . The  $(\alpha, \beta)$  are forward-backward messages. Their updating equations are derived in the appendix.

To determine the number of controller nodes  $\{|\mathcal{Z}_n|\}_{n=1}^N$ , the occupancy of a node is computed by checking if there is a positive reward assigned to it. For example, for action  $a$  and node  $i$ ,  $\hat{\rho}_{n,i}^a - \rho_{n,i}^a$  is the reward being assigned. If this quantity is greater than zero, then node  $i$  is visited. Summing over all actions gives the value of node  $i$ . Hence  $|\mathcal{Z}_n|$  can be computed based on the following formula

$$|\mathcal{Z}_n| = \sum_{i=1}^{\infty} \mathbb{I}\left(\sum_{a=1}^{|\mathcal{A}_n|} (\hat{\rho}_{n,i}^a - \rho_{n,i}^a) > 0\right). \quad (13)$$

The complete algorithm is described in Algorithm 1. Upon the convergence of Algorithm 1, point estimates of the decentralized policies may be obtained by calculating the expectation:  $\mathbb{E}[\hat{\mu}_n^i]$ ,  $\mathbb{E}[\hat{\pi}_{n,i}^a]$ , and  $\mathbb{E}[\hat{W}_{n,a,o}^{i,j}]$  (see the appendix for details).

**Table 1:** Computational Complexity of Algorithm 1.

| VAR       | BEST CASE  | WORST CASE  |
|-----------|--|---|
| $\alpha$  | $\Omega(N \mathcal{Z} ^2KT)$   | $O(N \mathcal{Z} ^2KT)$   |
| $\beta$   | $\Omega(N \mathcal{Z} ^2KT)$   | $O(N \mathcal{Z} ^2KT^2)$   |
| $\nu_t^k$ | $\Omega(K)$  | $O(KT)$   |
| $\Theta$  | $\Omega(N \mathcal{Z} ^2KT)$   | $O(N \mathcal{Z} ^2KT^2)$   |
| LB        | $\Omega( \mathcal{Z} ^2 \sum_{n=1}^N  \mathcal{A}_n   \mathcal{O}_n )$ | $O( \mathcal{Z} ^2 \sum_{n=1}^N  \mathcal{A}_n   \mathcal{O}_n )$ |

### 3.3 Computational complexity

The time complexity of Algorithm 1 for each iteration is summarized in Table 1, assuming the length of an episode is on the order of magnitude of  $T$ , and the number of nodes per controller is on the order of magnitude of  $|\mathcal{Z}|$ . In Table 1, the worst case refers to when there is a nonzero reward at every time step of an episode (dense rewards), while the best case is when nonzero reward is received only at the terminal step. Hence in general, the algorithm scales linearly with the number of episodes and the number of agents. The time dependency on  $T$  is between linear and quadratic. In any case, the computational complexity of Algorithm 1 is independent of the number of states, making it is scalable to large problems.

### 3.4 Exploration and Exploitation Tradeoff

Algorithm 1 assumes off-policy batch learning where trajectories are collected using a separate behavior policy. This is appropriate when data has been generated from real-world or simulated experiences without any input from the learning algorithm (e.g., learning from demonstration). Off-policy learning is efficient if the behavior policy is close to optimal, as in the case when expert information is available to guide the agents. With a random behavior policy, it may take a long time



for the policy to converge to optimality; in this case, the agents may want to exploit the policies learned so far to speed up the learning process.

An important issue concerns keeping a proper balance between exploration and exploitation to prevent premature convergence to a suboptimal policy, but allow the algorithm to learn quickly. Since the execution of Dec-POMDP policies is decentralized, it is difficult to design an efficient exploration strategy that guarantees optimality. [35] count the visiting frequency of FSC nodes and apply upper-confidence-bound style heuristic to select next controller nodes, and use  $\epsilon$ -greedy strategy to select actions. However  $\epsilon$ -greedy might be sample inefficient. [6] proposed a distributed learning approach where agents take turns to learn the best response to each other’s policies. This framework applies an R-max type of heuristic, using the counts of trajectories to distinguish known and unknown histories, to tradeoff exploration and exploitation. However, this method is confined to tree-based policies in finite-horizon problems, and requires synchronized multi-agent learning.

To better accommodate our Bayesian policy learning framework for RL in infinite-horizon Dec-POMDPs, we define an auxiliary FSC,  $\Omega_n = \langle \mathcal{Y}, \mathcal{O}_n, \mathcal{Z}_n, W_n, \mu_n, \varphi_n \rangle$ , to represent the policy of each agent in balancing exploration and exploitation. To avoid confusion, we refer to  $\Theta_n$  as a primary FSC. The only two components distinguishing  $\Psi_n$  from  $\Theta_n$  are  $\mathcal{Y}$  and  $\varphi_n$ , where  $\mathcal{Y} = \{0, 1\}$  encodes exploration ( $y = 1$ ) or exploitation ( $y = 0$ ), and  $\varphi_n = \{\varphi_y^{n,z}\}$  with  $\varphi_y^{n,z}$  denoting the probability of agent  $n$  choosing  $y$  in  $z$ . One can express  $p(y_{n,t}|h_{n,t}, \Psi_n)$  in the same way as one expresses  $p(a_{n,t}|h_{n,t}, \Theta_n)$  (which is described in section 2.1). The behavior policy  $\Pi_n$  of agent  $n$  is given as

$$p^{\Pi_n}(a|h, \Theta_n, \Omega_n) = \sum_{y=0,1} p(a|y, h)p(y|h, \Omega_n), \quad (14)$$

where  $p(a|y = 0, h) \equiv p(a|h, \Theta_n)$  is the primary FSC policy, and  $p(a|y = 1, h)$  is the exploration policy of agent  $n$ , which is usually a uniform distribution.

The behavior policy in (14) has achieved significant success in the single-agent case [13, 21, 22]. Here we extend it to the multi-agent case (centralized learning and decentralized exploration/execution) and provide empirical evaluation in the next section.

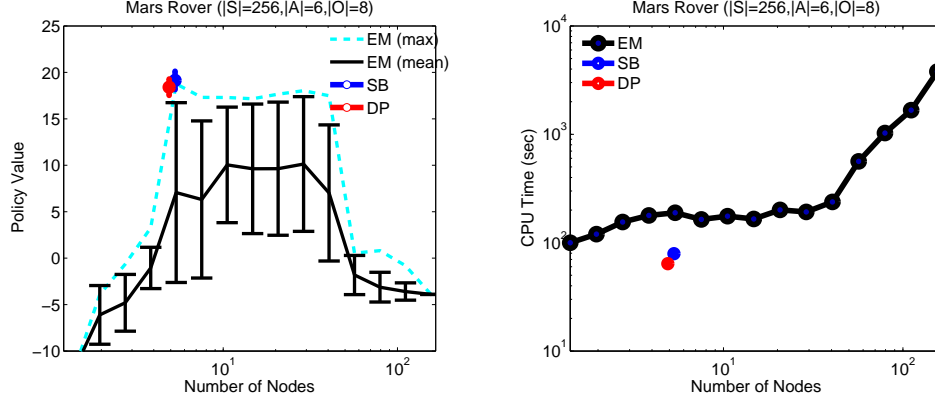
## 4 Experiments

The performance of the proposed algorithms are evaluated on five benchmark problems [1] and a large-scale problem (traffic control) [35]. The experimental procedure in [35] was used for all the results reported here. For Dec-SBPR, the hyperparameters in (8) are set to  $c = 0.1$  and  $d = 10^{-6}$  to promote sparse usage of FSC nodes.<sup>4</sup> The policies are initialized as FSCs converted from the episodes with the highest rewards using a method similar to [5].

**Learning variable-size FSC vs learning fixed-size FSC** To demonstrate the advantage of learning variable-size FSCs, Dec-SBPR is compared with an implementation of the previous EM algorithm [35]. The comparison is for the Mars Rover problem using  $K = 300$  episodes<sup>5</sup> to learn the FSCs and evaluating the policy by the discounted accumulated reward averaged over 100 test episodes of 1000 steps. Here, we consider off-policy learning and apply a semi-random policy to

<sup>4</sup>These values were chosen for testing, but our approach is robust to other values of  $c$  and  $d$ .

<sup>5</sup>Using smaller training sample size  $K$ , our method can still perform robustly, as it is shown in the appendix.

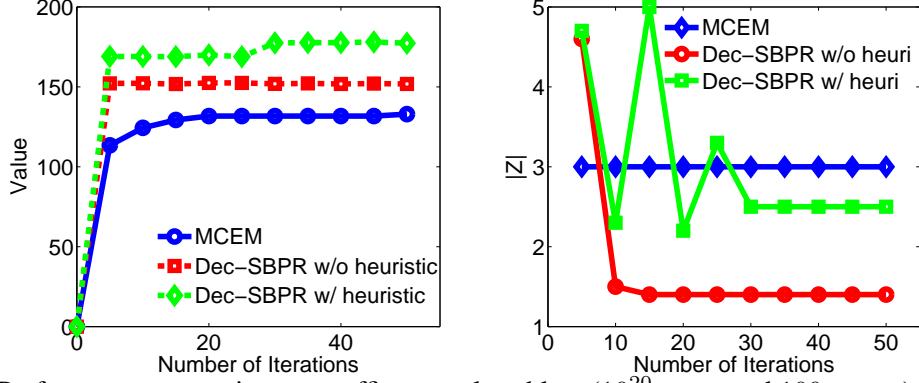


**Figure 1:** Comparison between the variable-size controller learned by Dec-SBPR and fixed-size controllers learned by EM. Left: testing value; Right: averaged computation time. Although EM has less computational complexity than our VB algorithm per iteration, empirically, the VB algorithm uses less time to reach convergence. Moreover, the average value of using the SB prior and its special Dirichlet instance (DP) are close to the best result of EM (dotted sky-blue line), but are much better than the average results of EM (solid back line). Using the SB prior achieves slightly better performance than the DP, which can be explained the flexibility of SB prior, as explained by the last paragraph of section 3.1.

**Table 2:** Performance of Dec-SBPR on benchmark problems compared to other state-of-art algorithms. Shows policy values (higher value indicates better performance) and CPU times of all algorithms, and the average controller size  $|Z|$  inferred by Dec-SBPR.

| PROBLEMS ( $ S ,  A ,  O $ ) | POLICY LEARNING (UNKNOWN MODEL) |       |       |                      |       |      | PLANNING (KNOWN MODEL) |               |       |         |               |         |
|------------------------------|---------------------------------|-------|-------|----------------------|-------|------|------------------------|---------------|-------|---------|---------------|---------|
|                              | DEC-SBPR(FIXED ITERATION)       |       |       | DEC-SBPR(FIXED TIME) |       |      | PERIEM                 |               |       | FB-HSVI |               |         |
|                              | VALUE                           | $ Z $ | TIME  | VALUE                | $ Z $ | TIME | VALUE                  | $ Z $         | TIME  | VALUE   | $ Z $         | TIME    |
| DEC-TIGER (2, 3, 3)          | -18.63                          | 6     | 96s   | -19.42               | 8     | 20s  | -32.31                 | 3             | 20s   | 9.42    | $7 \times 10$ | 6540s   |
| BROADCAST (4, 2, 5)          | 9.20                            | 2     | 7s    | 9.27                 | 2     | 24s  | 9.15                   | 3             | 24s   | 13.45   | 52            | 6.0s    |
| RECYCLING ROBOTS (3, 3, 2)   | 31.26                           | 3     | 147s  | 25.16                | 2     | 19s  | 30.78                  | 3             | 19s   | 9.27    | 102           | 19.8s   |
| BOX PUSHING (100, 4, 5)      | 77.65                           | 14    | 290s  | 58.27                | 9     | 32s  | 59.95                  | 3             | 32s   | 31.93   | 108           | 0s      |
| MARS ROVERS (256, 6, 8)      | 20.62                           | 5     | 1286s | 15.2                 | 6     | 160s | 8.16                   | 3             | 160s  | 224.43  | 331           | 1715.1s |
|                              |                                 |       |       |                      |       |      | 18.13                  | $3 \times 10$ | 7132s | 26.94   | 136           | 74.31s  |

collect samples. Specifically, the learning agent is allowed access to episodes collected by taking actions according to a POMDP algorithm (point-based value iteration (PBVI) [29]). Let  $\epsilon$  be the probability that the agents follow the PBVI policy and  $1 - \epsilon$  be the probability that the agents take random actions. This procedure mimics the approach used in previous work [35]. The results with  $\eta = 0.3$  are reported in Figure 1, which shows the exact value and computation time as a function of the number of controller nodes  $|Z|$ . As expected, for the EM algorithm, when  $|Z|$  is too small, the FSCs cannot represent the optimal policy (under-fitting), and when the number of nodes is too large, FSCs overfits a limited amount of data and perform poorly. Even if  $|Z|$  is set to the number inferred by Dec-SBPR, EM can still suffer severely from initialization and local maxima issues, as can be seen from a large error-bar. By setting a high truncation level ( $|Z| = 50$ ), Dec-SBPR employs Algorithm 1 to integrate out the uncertainty of the policy representation (under the SB prior). As a result, Dec-SBPR can infer both the number of nodes that is needed ( $\approx 5$ ) and optimal controller parameters simultaneously. Furthermore, this inference is done with less computation time and with a higher value and improved robustness (low variance of test value) than EM.



**Figure 2:** Performance comparison on traffic control problem ( $10^{20}$  states and 100 agents). Test reward (left) and inferred controller size (right) of Dec-SBPR, as a function of algorithmic iteration.

**Comparison with other methods** The performance of Dec-SBPR is also compared to several state-of-art methods, including: Monte Carlo EM (MCEM) [35]. Similar to Dec-SBPR, MCEM is a policy-based RL approach. We apply the exploration-exploitation strategy described in section 3.4 and follow the same experimental procedure in [35] to report the results<sup>6</sup>. The rewards after running a fixed number of iterations and a fixed amount of time are summarized (respectively) in Table 2 (the first column under policy-learning category). Dec-SBPR is shown to achieve better policy values than MCEM on all problems<sup>7</sup>. These results can be explained by the fact that EM is (more) sensitive to initialization and (more) prone to local optima. Moreover, by fixing the size of the controllers, the optimal policy from EM algorithms might be over/under fitted. By using a Bayesian nonparametric prior, Dec-SBPR learns the policy with variable-size controllers, allowing more flexibility for representing the optimal policy. We also show the result of Dec-SBPR running the same amount of clock time as MCEM (Dec-SBPR (fixed time)), which indicates Dec-SBPR can achieve a better trade-off between policy value and learning time than MCEM.

Finally, Dec-SBPR is compared to Periodic EM (PeriEM) [27] and FB-HSVI [15], two state-of-art planing methods (with known models) for generating controllers. Because having a Dec-POMDP model allows more accurate value function calculations than a finite number of trajectories, the value of PeriEM and FB-HSVI are treated as upper-bounds for the policy-based methods. Our Dec-SBPR approach can sometimes outperform PeriEM, but produces lower value than FB-HSVI. FB-HSVI is a boundedly-optimal method, showing that Dec-SBPR can produce near optimal solutions in some of these problems and produces solutions that are much closer to the optimal than previous RL methods. It is also worth noting that neither PeriEM nor FB-HSVI can scale to large problems (such as the one discussed below), while by using a policy-based RL approach, Dec-SBPR can scale well.

**Scaling up to larger domains** To demonstrate scalability to both large problem sizes and large numbers of agents, we test our algorithm on a traffic problem [35], with  $10^{20}$  states. Here, there are 100 agents controlling the traffic flow at  $10 \times 10$  intersections with one agent located at each intersection. Except for MCEM, no previous Dec-POMDPs algorithms are able to solve such large

<sup>6</sup>The learning curves of Dec-SBPR are shown in the appendix.

<sup>7</sup>The results are provided by personal communication with its authors and run on the same benchmarks that are available online.

problems.

Since the authors in [35] use a hand-coded policy (comparing the traffic flow between two directions) as a heuristic for generating training trajectories, we also use such a heuristic for a fair comparison. In addition, to examine the effectiveness of the exploration-exploitation strategy described in Section 3.4, we also consider the case where the initial behavior policy is random and then it is optimized as discussed. From Figure 2, we can see that, with the help of the heuristic, Dec-SBPR can achieve the best performance. Without using the heuristic (by just using our exploration-exploitation strategy), in a few iterations, Dec-SBPR is able to produce a higher quality policy than MCEM. Moreover, the inferred number of FSC nodes (averaged over all agents) is smaller than the number preselected by MCEM. This shows that not only can Dec-SBPR scale to large problems, but it can also produce higher-quality solutions than other methods for those large problems.

## 5 Conclusions

The paper presented a scalable Bayesian nonparametric policy representation and an associated learning framework (Dec-SBPR) for generating decentralized policies in Dec-POMDPs. An new exploration-exploitation method, which extends the popular  $\epsilon$ -greedy method, was also provided for reinforcement learning in Dec-POMDPs. Experimental results show Dec-SBPR produces higher-quality solutions than the state-of-art policy-based method, and has the additional benefit of inferring the number of nodes needed to represent the optimal policy. The resulting method is also scalable to large domains (in terms of both the number of agents and the problem size), allowing high-quality policies for large Dec-POMDPs to be learned efficiently from data.

## Acknowledgments

Acknowledgments This research was supported by the US Office of Naval Research (ONR) under MURI program award #N000141110688 and NSF award #1463945.

## Appendices

### A Proof of Theorem 3:(Mean-field) Variational Bayesian (VB) Inference for DEC-SBPR

Under the standard variational theory [7, 11], minimizing the KL divergence between  $q(\Theta, z)$  and  $p(\Theta, z|\mathcal{D})$  is equivalent to minimizing the lower bound of log marginal likelihood (empirical value function for our case). Using Jensen’s inequality, we can obtain the following lower bound of the log marginal value function

$$\ln \widehat{V}(\mathcal{D}^{(K)}) = \ln \frac{1}{K} \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{z_0^k, \dots, z_t^k} \int q_t^k(z_{0:t}^k) q(\Theta) q(\eta) \frac{\tilde{r}_t^k p(\Theta) p(\eta) p(\vec{a}_{0:t}^k, \vec{z}_{0:t}^k | \vec{o}_{1:t}^k, \Theta)}{q_t^k(\vec{z}_{0:t}^k) q(\Theta) q(\eta)} d\Theta d\eta$$

$$\begin{aligned}
&\geq \frac{1}{K} \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\vec{z}_0^k, \dots, \vec{z}_t^k} \int q_t^k(\vec{z}_{0:t}^k) q(\Theta) q(\eta) \ln \frac{\tilde{r}_t^k p(\Theta) p(\eta) p(\rho) p(\vec{a}_{0:t}^k | \vec{o}_{1:t}^k, \Theta)}{q_t^k(\vec{z}_{0:t}^k) q(\Theta) q(\eta)} d\Theta d\eta \\
&= \ln \widehat{V}(\mathcal{D}^{(K)}) - KL \left( \left\{ q_t^k(\vec{z}_{0:t}^k) q(\Theta) q(\eta) \right\} \middle| \middle| \frac{\xi_t^k}{\widehat{V}(\mathcal{D}^{(K)})} p(\vec{z}_{0:t}^k, \Theta, \eta | \vec{a}_{0:t}^k, \vec{o}_{1:t}^k) \right) \\
&\stackrel{def}{=} \text{LB} \left( \left\{ \left\{ q_t^k(z_{n,0:t}^k) \right\}_{k,t}, q(\Theta_n), q(\eta_n) \right\}_{n=1 \dots N} \right) \tag{15}
\end{aligned}$$

where  $\xi_t^k = \tilde{r}_t^k \int p(\vec{a}_{0:t}^k | \vec{o}_{1:t}^k, \Theta) p(\Theta) d\Theta$  and  $\tilde{r}_t^k = \gamma^t r_t^k / p^\Psi(\vec{a} | \vec{h}_t^k)$ . We assume  $p(\Theta) = \prod_{n=1}^N p(\Theta)$  and  $p(\vec{a}_{0:t}, \vec{z}_{0:t} | \vec{o}_{0:t}, \Theta) = \prod_{n=1}^N p(a_{n,0:t}, z_{n,0:t} | o_{n,0:t}, \Theta_n)$  to accommodate decentralized policy representations.

To derive the VB updating equations, we rewrite the lower bound in equation (15) as follows

$$\begin{aligned}
&\text{LB} \left( \left\{ \left\{ q_t^k(z_{n,0:t}^k) \right\}_{k,t}, q(\Theta_n), q(\eta_n) \right\}_{n=1, \dots, N} \right) \\
&= - \int q(\Theta) \ln \frac{q(\Theta)}{p(\Theta)} d\Theta - \int q(\eta) \ln \frac{q(\eta)}{p(\eta)} d\eta - \frac{1}{K} \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\vec{z}_0^k, \dots, \vec{z}_t^k} q_t^k(\vec{z}_{0:t}^k) \ln (q_t^k(\vec{z}_{0:t}^k)) \\
&\quad + \frac{1}{K} \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\vec{z}_0^k, \dots, \vec{z}_t^k} \prod_{n=1}^N \int q_t^k(\vec{z}_{0:t}^k) q(\Theta_n) q(\eta_n) \ln \left( \tilde{r}_t^k \prod_{n=1}^N p_t^k(a_{n,0:t}^k, z_{n,0:t}^k | o_{n,1:t}^k, \Theta_n) \right) d\Theta_n d\eta_n \tag{16}
\end{aligned}$$

The VB Inference algorithm for DEC-SBPR is based on maximizing LB w.r.t. the distribution of the joint DEC-SBPR parameters  $\left\{ \left\{ q_t^k(z_{n,0:t}^k) \right\}_{k,t}, q(\Theta_n), q(\eta_n) \right\}_{n=1, \dots, N}$ , which can be achieved by alternating the following steps.

**Update the distribution of nodes (VB E-step):** Keeping  $\{q(\Theta_n)\}_{n=1, \dots, N}$  and  $\{q(\eta_n)\}_{n=1, \dots, N}$  fixed, solve  $\max_{\{q_t^k(z_{n,0:t}^k)\}} \text{LB} \left( \left\{ \left\{ q_t^k(z_{n,0:t}^k) \right\}_{k,t}, q(\Theta_n), q(\eta_n) \right\}_{n=1, \dots, N} \right)$ ,  $\forall n, k, t$  subject to the normalization constraint for  $q_t^k(z_{n,0:t}^k)$ . In this step, we construct the Lagrangian

$$F_{q_t^k(z_{n,0:t}^k)} = \text{LB} \left( \left\{ \left\{ q_t^k(z_{n,0:t}^k) \right\}_{k,t}, q(\Theta_n), q(\eta_n) \right\}_{n=1, \dots, N} \right) - \lambda \left( K - \sum_{k,t,z_{n,0:t}^k} \prod_{n=1}^N q_t^k(z_{n,0:t}^k) \right), \tag{17}$$

then take derivative w.r.t  $q_t^k(z_{n,0:t}^k)$  and set the result to zero

$$\begin{aligned}
\frac{\partial F_{q_t^k(z_{n,0:t}^k)}}{\partial (q_t^k(z_{n,0:t}^k))} &= \frac{1}{K} \int p(\Theta_n) \ln \tilde{r}_t^k p(a_{n,0:t}^k, z_{n,0:t}^k | o_{1:t}^k, \Theta_n) d\Theta_n \\
&\quad + \frac{1}{K} \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\vec{z}_0^k, \dots, \vec{z}_t^k} \prod_{i \neq n} \int q_t^k(z_{i,0:t}^k) q(\Theta_n) q(\eta_n) \ln \left( \tilde{r}_t^k \prod_{n=1}^N p_t^k(a_{n,0:t}^k, z_{n,0:t}^k | o_{n,1:t}^k, \Theta_n) \right) d\Theta_n d\eta_n \\
&\quad - \frac{1}{K} \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\vec{z}_0^k, \dots, \vec{z}_t^k} \prod_{i \neq n} q_t^k(z_{i,0:t}^k) \ln \left( \prod_{n=1}^N q_t^k(z_{n,0:t}^k) \right) - \frac{1}{K} \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\vec{z}_0^k, \dots, \vec{z}_t^k} \prod_{i \neq n} q_t^k(z_{i,0:t}^k)
\end{aligned}$$

$$+\lambda \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{z_0^k, \dots, z_t^k} \prod_{i \neq n} q_t^k(z_{i,0:t}^k) = 0, \quad (18)$$

which is solved to give the distribution of nodes  $z_{n,0:t}^k$  for the  $n^{th}$  agent

$$\begin{aligned} q_t^k(z_{n,0:t}^k) &= \frac{\tilde{r}_t^k}{C_z} \exp \left\{ \int q(\Theta_n) \ln p(a_{n,0:t}^k, z_{n,0:t}^k | o_{n,1:t}^k, \Theta_n) d\Theta_n \right\} \\ &= \frac{\tilde{r}_t^k}{C_z} \exp \left\{ \ln \mu_n^{z_{n,0}^k} + \sum_{\tau=0}^t \langle \ln \pi_{n,z_{n,\tau}^k}^{a_{n,\tau}^k} \rangle_{p(\pi|\hat{\rho})} + \sum_{\tau=1}^t \langle \ln W_{n,a_{n,\tau-1}^k, o_{\tau}^k}^{z_{n,\tau-1}^k, z_{n,\tau}^k} \rangle_{p(W|\hat{\sigma}, \hat{\eta})} \right\} \\ &= \frac{\tilde{r}_t^k}{C_z} \mu_n^{z_{n,0}^k} \tilde{\pi}_{n,z_{n,0}^k}^{a_{n,0}^k} \prod_{\tau=1}^t \tilde{\pi}_{n,a_{n,\tau-1}^k, o_{\tau}^k}^{z_{n,\tau-1}^k, z_{n,\tau}^k} \end{aligned} \quad (19)$$

where

$$\begin{aligned} \tilde{\pi}_a^{n,i} &= \exp \left\{ \langle \ln \pi_a^{n,i} \rangle_{p(\pi|\hat{\rho})} \right\} = \exp \left\{ \langle \psi(\hat{\rho}_a^{n,i}) - \psi(\sum_{m=1}^{|\mathcal{A}|} \hat{\rho}_m^{n,i}) \rangle \right\} \\ \widetilde{W}_{n,a,o}^{i,1} &= \exp \left\{ \langle \ln W_{n,a,o}^{i,1} \rangle_{p(W|\hat{\sigma}, \hat{\eta})} \right\} = \exp \left\{ \langle \ln V_{n,a,o}^{i,1} \rangle_{p(V|\hat{\sigma}, \hat{\eta})} \right\}, \\ \widetilde{W}_{n,a,o}^{i,j} &= \exp \left\{ \langle \ln W_{n,a,o}^{i,j} \rangle_{p(W|\hat{\sigma}, \hat{\eta})} \right\} = \exp \left\{ \langle \ln V_{n,a,o}^{i,j} \rangle_{p(V|\hat{\sigma}, \hat{\eta})} + \sum_{m=1}^{j-1} \langle \ln(1 - V_{n,a,o}^{i,m}) \rangle_{p(V|\hat{\sigma}, \hat{\eta})} \right\}, \\ &\quad \text{for } j=2, \dots, |\mathcal{Z}_n|-1 \\ \widetilde{W}_{n,a,o}^{i,|\mathcal{Z}_n|} &= \exp \left\{ \langle \ln W_{n,a,o}^{i,|\mathcal{Z}_n|} \rangle_{p(W|\hat{\sigma}, \hat{\eta})} \right\} = \exp \left\{ \sum_{m=1}^{|\mathcal{Z}_n|} \langle \ln(1 - V_{n,a,o}^{i,m}) \rangle_{p(V|\hat{\sigma}, \hat{\eta})} \right\} \end{aligned}$$

with

$$\begin{aligned} \langle \ln V_{n,a,o}^{i,j} \rangle_{p(V|\hat{\sigma}, \hat{\eta})} &= \psi(\sigma_{n,a,o}^{i,j} + \omega_{n,a,o}^{i,j}) - \psi\left(\sigma_{n,a,o}^{i,j} + \langle \eta_{n,a,o}^{i,j} \rangle + \sum_{l=j}^{|\mathcal{Z}_n|} \omega_{n,a,o}^{i,l}\right), \\ \langle \ln(1 - V_{n,a,o}^{i,j}) \rangle_{p(V|\hat{\sigma}, \hat{\eta})} &= \psi\left(\langle \Sigma_{n,a,o}^{i,j} \rangle + \sum_{l=j+1}^{|\mathcal{Z}_n|} \omega_{n,a,o}^{i,l}\right) - \psi\left(\sigma_{n,a,o}^{i,j} + \langle \Sigma_{n,a,o}^{i,j} \rangle + \sum_{l=j}^{|\mathcal{Z}_n|} \omega_{n,a,o}^{i,l}\right), \\ \langle \eta_{n,a,o}^{i,j} \rangle_{p(\eta|\hat{c}, \hat{d})} &= \hat{c}_{n,a,o}^{i,j} / \hat{d}_{n,a,o}^{i,j} \end{aligned}$$

where  $\psi(\cdot)$  is digamma function,  $\omega_{n,a,o}^{i,j}$  is the reward (soft-count) allocated to the transition from state  $i$  to  $j$ , both  $\hat{c}_{n,a,o}^{i,j}$  and  $\hat{d}_{n,a,o}^{i,j}$  are the posterior parameters of  $\eta_{n,a,o}^{i,j}$ .

**Update the sufficient statistics (VB-M Step)** In VB-M Step, the distribution of nodes  $\{q_t^k(z_{0:t}^k)\}$ ,  $\forall k, t$  are fixed, the objective is to solve  $\max_{q(\Theta_n), q(\eta_n)} \text{LB} \left( \left\{ \{q_t^k(z_{n,0:t}^k)\}_{k,t}, q(\Theta), q(\eta_n) \right\}_{n=1, \dots, N} \right)$ ,  $\forall n$  subject to the normalization constraint that  $\int q(\Theta) d\Theta = 1$ . First we consider finding  $q(\Theta_n)$ . To that end, we construct the Lagrangian

$$F_{q_{\Theta_n}} = \text{LB} \left( \left\{ \{q_t^k(z_{n,0:t}^k)\}_{k,t}, q(\Theta_n), q(\eta_n) \right\}_{n=1, \dots, N} \right) - \lambda \left( 1 - \int q(\Theta) d\Theta \right). \quad (20)$$

Then

$$\frac{\partial F_{q_\Theta}}{\partial(q(\Theta))} = \frac{1}{K} \sum_{k,t,z^k_{n,0:t}} \int \prod_{n=1}^N q(\eta_n) q_t^k(z^k_{n,0:t}) \ln \tilde{r}_t^k \prod_{n=1}^N p(\eta_n) p(a^k_{n,0:t}, z^k_{n,0:t} | o^k_{1:t}, \Theta_n) d\eta$$

$$-1 - \ln \frac{q(\Theta_n)}{p(\Theta_n)} + \lambda = 0$$

$\Rightarrow$

$$\begin{aligned} q(\Theta) &= \prod_{n=1}^N \frac{p(\Theta_n)}{e^{1-\lambda}} \exp \left\{ \frac{1}{K} \sum_{k,t,z} q_t^k(z^k_{n,0:t}) \ln p(a^k_{n,0:t}, z^k_{n,0:t} | o^k_{1:t}, \Theta_n) \right\} \\ &= \prod_{n=1}^N \frac{p(\Theta_n)}{e^{1-\lambda}} \exp \left\{ \frac{1}{K} \sum_{k,t} K \nu_t^k \left[ \sum_{\tau=1}^t \sum_{i=1}^{|\mathcal{Z}_n|} \phi_{t,\tau}^{n,k}(i) \ln \pi_i^{n,a^k_{n,\tau}} + \sum_{\tau=1}^t \sum_{i,j=1}^{|\mathcal{Z}_n|} \xi_{t,\tau}^{n,k}(i,j) \ln W_{i,j}^{n,a^k_{\tau-1}, o^k_\tau} \right] \right\} \\ &= \prod_{n=1}^N \frac{p(\Theta_n)}{e^{1-\lambda}} \prod_{i=1}^{|\mathcal{Z}_n|} \prod_{a=1}^{|\mathcal{A}_n|} \left\{ \left\{ [\pi_{n,a}^i]^{\sum_{k,t,\tau} \nu_t^k \phi_{t,\tau}^k(i) \mathbb{I}(a^k_{n,\tau}, a)} \prod_{j=1}^{|\mathcal{Z}_n|} \prod_{o=1}^{|\mathcal{O}_n|} [W_{n,a,o}^{i,j}]^{\sum_{k,t,\tau} \nu_t^k \xi_{t,\tau-1}^{n,k}(i,j) \mathbb{I}(a^k_{n,\tau-1}, a) \mathbb{I}(o^k_{n,\tau}, o)} \right\} \right\} \\ &\quad \text{(from the relation between stick-breaking weights } \mathbf{p} \text{ and independent beta random variables in } \mathbf{V} \text{ (28))} \\ &= \prod_{n=1}^N \frac{p(\Theta_n)}{e^{1-\lambda}} \prod_{i=1}^{|\mathcal{Z}_n|} \prod_{a=1}^{|\mathcal{A}_n|} \left\{ [\pi_{n,i}^a]^{\sum_{k,t,\tau} \nu_t^k \phi_{t,\tau}^{n,k}(i) \mathbb{I}(a, a^{n,k}_\tau)} \right. \\ &\quad \times \left. \prod_{j=1}^{|\mathcal{Z}_n|} \prod_{o=1}^{|\mathcal{O}_n|} [V_{n,a,o}^{i,j}]^{\sum_{k,t,\tau} \nu_t^k \xi_{t,\tau-1}^{n,k}(i,j) \mathbb{I}(a^k_{n,\tau-1}, a) \mathbb{I}(o^k_{n,\tau}, o)} \right\} \\ &= \prod_{n=1}^N \text{Dir}(\hat{\rho}_{n,i}) \prod_{a=1}^{|\mathcal{A}_n|} \prod_{o=1}^{|\mathcal{O}_n|} \text{GDD}(\hat{\sigma}_{n,a,o}^{i,\cdot}, \hat{\eta}_{n,a,o}^{i,\cdot}) \end{aligned} \quad (21)$$

where

$$\begin{aligned} \hat{\rho}_{n,i}^a &= \rho_{n,i}^a + \frac{1}{K} \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\tau=1}^t \nu_t^k \phi_{t,\tau-1}^{n,k}(i) \mathbb{I}(a^k_{\tau-1} = a) \\ \hat{\sigma}_{n,a,o}^{i,j} &= \sigma_{n,a,o}^{i,j} + \omega_{n,a,o}^{i,j} \\ &= \sigma_{n,a,o}^{i,j} + \frac{1}{K} \sum_{k=1}^K \sum_{t=0}^{T_k} \sum_{\tau=1}^t \nu_t^k \xi_{t,\tau-1}^{n,k}(i,j) \mathbb{I}(a^k_{\tau-1} = a) \mathbb{I}(o^k_\tau = o) \\ \hat{\eta}_{n,a,o}^{i,j} &= \eta_{n,a,o}^{i,j} + \sum_{l=j+1}^{|\mathcal{Z}_n|} \omega_{n,a,o}^{i,l} \\ &= \eta_{n,a,o}^{i,j} + \frac{1}{K} \sum_{l=j+1}^{|\mathcal{Z}_n|} \sum_{k=1}^K \sum_{t=0}^{T_k} \hat{\nu}_t^k \sum_{\tau=1}^t \xi_{t,\tau-1}^{n,k}(i,j) \mathbb{I}(a^k_{\tau-1} = a) \mathbb{I}(o^k_\tau = o), \end{aligned} \quad (22)$$

and  $\nu_t^k$  is the marginalized re-weighted reward computed by equation 6.

To find  $q(\eta)$ , we construct the Lagrange,

$$F_{q(\Sigma)} = \text{LB}(\{q_t^k\}, g(\Theta), q(\eta)) - \lambda \left( 1 - \int q(\eta) d\eta \right) \quad (23)$$

$$\frac{\partial F_\eta}{\partial(q(\eta))} = \frac{1}{K} \int q(\Theta) \ln p(\eta) p(\Theta) d\Theta - \frac{1}{K} \ln q(\eta) + \lambda = 0 \quad (24)$$

$\Rightarrow$

$$\begin{aligned} q(\eta) &\propto \exp \left\{ \int q(\Theta) \ln p(\eta) p(\Theta|\eta) d\Theta \right\} \\ &\propto \exp \left\{ \int q(\Theta) \ln p(\Theta|\eta) d\Theta + \ln p(\eta) \right\} \\ &\propto p(\eta) \exp \{ \langle \ln p(\Theta) \rangle \} \\ &= \prod_{n=1}^N \prod_{a=1}^{|\mathcal{A}_n|} \prod_{j=1}^{|\mathcal{O}_n|} \prod_{i=1}^{|\mathcal{Z}_n|} \prod_{j=1}^{|\mathcal{Z}_n|} \text{Ga}(\eta_{i,j}^{n,a,o}; c, d) \frac{\Gamma(\sigma_{i,j}^{n,a,o} + \eta_{i,j}^{n,a,o})}{\Gamma(\sigma_{i,j}^{n,a,o}) \Gamma(\eta_{i,j}^{n,a,o})} V_{i,j}^{n,a,o} \sigma_{i,j}^{n,a,o-1} (1 - V_j^{n,a,o})^{\eta_{i,j}^{n,a,o}-1} \\ &\approx \prod_{n=1}^N \prod_{a=1}^{|\mathcal{A}_n|} \prod_{j=1}^{|\mathcal{O}_n|} \prod_{i=1}^{|\mathcal{Z}_n|} \prod_{j=1}^{|\mathcal{Z}_n|} \frac{\Gamma(\sigma_{i,j}^{n,a,o} + \eta_{i,j}^{n,a,o})}{\Gamma(\eta_{i,j}^{n,a,o})} \{\eta_{i,j}^{n,a,o}\}^{c-1} \exp \{ -\eta_{i,j}^{n,a,o} (d - \ln(1 - V_j^{n,a,o})) \} \end{aligned} \quad (25)$$

One can set  $\sigma_{i,j}^{n,a,o} = 1$ , in this case, the VB approximation of  $q(\sigma)$  is a product of independent gamma distributions. However, when  $\sigma_{i,j}^{n,a,o} \neq 1$ , (24) is no longer a gamma distribution (the prior and likelihood are not conjugate). To solve this issue, one might consider the VB inference method for non-conjugate priors [33], by which we consider a point estimate of  $\eta$ , such as  $q(\eta)$  is maximized. One way to obtain the maximum estimate of  $\eta$  is to solve  $\frac{\partial q(\eta)}{\partial \eta} = 0$ , however this operation involves taking derivative w.r.t gamma functions, which does not have a simple form solution. To circumvent this difficult, we use grid search. To make the search more efficient, we use the bounds of  $\frac{\Gamma(a+x)}{\Gamma(x)}$  to give an initial estimate of the searching range. The bounds are from Wendel's double Inequality [30].

$$x(x+a)^{a-1} \leq \Gamma(a+x)/\Gamma(x) \leq x^a \quad (26)$$

where  $x, a > 0$ .

## B Some Basics of Stick-breaking Priors

Here we provide the definition of stick-breaking prior (SBP), its connection to generalized Dirichlet distribution and the corresponding posterior inference, as well as the main statistics characteristics which are useful for developing the inference methods in our paper. For more detailed mathematical treatment, readers are referred to [18, 34]

**Definition 4.** *The stick-breaking priors [18] are almost surely discrete random probability measures  $\mathcal{P}$  over the measurable space  $(\Omega, \mathcal{B})$  which are partitioned into  $d$  disjoint regions with  $\Omega = \cup \mathcal{B}_k$  for  $1, \dots, d$ . It is expressed as*

$$\mathcal{P}_{a,b}^d = \sum_{k=1}^d p_i \delta_{\Theta_i} \quad (27)$$

and

$$p_1 = V_1 \quad \text{and} \quad p_i = (1 - V_1)(1 - V_2) \cdots (1 - V_{i-1})V_i, i > 2 \quad (28)$$

are the weights with  $V_i$  are independent  $\text{Beta}(a_k, b_k)$  random variables for  $a_i, b_i > 0$ .



SBP allows Beta-distributed RVs  $V_i, \forall i$  and the atoms  $\Theta_i, \forall i$  associated with the resulting weights to be drawn simultaneously.

### B.1 Stick-breaking prior and generalized Dirichlet distribution

We denote  $\mathbf{p} \sim SB(\mathbf{v}, \mathbf{w})$  as constructing  $\mathbf{p}$  as an infinite process ( $d \rightarrow \infty$ ) as (28), and  $\mathbf{p} \sim GDD(\mathbf{v}, \mathbf{w})$  when  $\mathbf{p}$  is finite. Here, GDD stands for generalized Dirichlet distribution. To see the connection between SBP and GDD, set the truncation level (number of occupied states) to  $d$  with  $p_{d+1} = 1 - \sum_{i=1}^d p_i$ , then we can write down the density function of  $\mathbf{V} = (V_1, \dots, V_d)$  as

$$f(\mathbf{V}) = \prod_{i=1}^d f(V_i) = \prod_{i=1}^d \frac{\Gamma(v_i + w_i)}{\Gamma(v_i)\Gamma(w_i)} V_i^{v_i-1} (1 - V_i)^{w_i-1}. \quad (29)$$

By changing variables from  $\mathbf{V}$  to  $\mathbf{p}$  and using the relation between  $\mathbf{V}$  and  $\mathbf{p}$  as described by (28), the density of  $\mathbf{p}$  can be obtained as follows,

$$f(\mathbf{p}) = \left| \frac{\partial \mathbf{V}}{\partial \mathbf{p}} \right| f(\mathbf{V}) = \prod_{i=1}^d \left( \frac{\Gamma(v_i + w_i)}{\Gamma(v_i)\Gamma(w_i)} p_i^{v_i-1} \left( 1 - \sum_{j=1}^i p_j \right)^{w_i-(v_{i+1}+w_{i+1})} \right) p_{d+1}^{w_d-1} \quad (30)$$

which has a mean and variance for an element  $p_i$ ,

$$\begin{aligned} \mathbb{E}[p_i] &= \frac{v_i \prod_{l=1}^{i-1} w_l}{\prod_{l=1}^i (v_l + w_l)}, \\ \mathbb{V}[p_i] &= \frac{v_i(v_i + 1) \prod_{l=1}^{i-1} w_l(w_l - 1)}{\prod_{l=1}^i (v_l + w_l)(v_l + w_l + 1)}. \end{aligned} \quad (31)$$

When  $w_i = \sum_{j=i+1}^K v_j$  for  $i < d$ , and keeping  $w_d = w_d$ , the GDD is equivalent to the standard Dirichlet distribution.

As a concrete example, consider the case  $d = 3$ , we have

$$\begin{aligned} p_1 &= V_1 & V_1 &= p_1 \\ p_2 &= (1 - V_1)V_2 & V_2 &= \frac{p_2}{1 - p_1} \\ p_3 &= (1 - V_1)(1 - V_2)V_3 & V_3 &= \frac{p_3}{1 - p_1 - p_2} \\ p_4 &= (1 - V_1)(1 - V_2)(1 - V_3)V_4 & V_4 &= \frac{p_4}{1 - p_1 - p_2 - p_3} = 1. \\ &= (1 - V_1)(1 - V_2)(1 - V_3) \end{aligned} \quad \Longleftrightarrow$$

Plugging these relations into (30), we can obtain

$$\begin{aligned} f(\mathbf{p}) &= \left| \frac{\partial \mathbf{V}}{\partial \mathbf{p}} \right| f(\mathbf{V}) \\ &= \prod_{i=1}^4 \left( 1 - \sum_{j=1}^{i-1} p_j \right) \frac{\Gamma(v_i + w_i)}{\Gamma(v_i)\Gamma(w_i)} \times p_i^{v_i-1} (1 - p_1)^{w_1-1} \cdot \left( \frac{p_2}{1 - p_1} \right)^{v_2-1} \left( 1 - \frac{p_2}{1 - p_1} \right)^{w_2-1} \\ &= \left( \frac{p_3}{1 - p_1 - p_2} \right)^{v_3-1} \left( 1 - \frac{p_2}{1 - p_1 - p_2} \right)^{w_3-1} \cdot \left( \frac{p_4}{1 - p_1 - p_2 - p_3} \right)^{v_4-1} \left( 1 - \frac{p_4}{1 - p_1 - p_2 - p_3} \right)^{w_4-1} \\ &= \prod_{i=1}^4 \frac{\Gamma(v_i + w_i)}{\Gamma(v_i)\Gamma(w_i)} p_i^{v_i-1} \left( 1 - \sum_{j=1}^i p_j \right)^{w_i-(v_{i+1}+w_{i+1})}. \end{aligned} \quad (32)$$

## B.2 Bayesian Inference for GDD

Given a set of discrete observations  $\{X_n\} \stackrel{i.i.d.}{\sim} \text{Discrete}(\mathbf{p})$ , and the prior  $\mathbf{p} \sim SB(v, w)$ , the posterior of  $\mathbf{p}$  can be written down as

$$\begin{aligned}
p(\mathbf{p}|\{X_n\}) &\propto \prod_{n=1}^N \prod_{i=1}^d p_i^{\mathbb{I}(X_n=i)} \prod_{i=1}^d \left( \frac{\Gamma(v_i + w_i)}{\Gamma(v_i)\Gamma(w_i)} p_i^{v_i-1} (1 - \sum_{j=1}^i p_j)^{w_i-(v_{i+1}+w_{i+1})} \right) p_{d+1}^{w_d-1} \\
&\propto \prod_{i=1}^d (V_i \prod_{j=1}^{i-1} (1 - V_j))^{ \sum_{n=1}^N \mathbb{I}(X_n=i) } V_i^{v_i-1} (1 - V_i)^{w_i-1} \\
&\propto \prod_{i=1}^d V_i^{v_i + \sum_{n=1}^N \mathbb{I}(X_n=i)-1} (1 - V_i)^{w_i + \sum_{j>i} \sum_{n=1}^N \mathbb{I}(X_n=i)-1} \\
&= GDD(\mathbf{v}', \mathbf{w}')
\end{aligned} \tag{33}$$

where the posterior hyper-parameters are updated as  $v'_i = v_i + \sum_{n=1}^N \mathbb{I}(X_n = i)$  and  $w'_i = w_i + \sum_{j>i} \sum_{n=1}^N \mathbb{I}(X_n = i)$ , where  $\mathbb{I}(\cdot)$  is an indicator function with value equal to one when the argument is true and zero otherwise.

## C The Computation of Forward and Backward Variables $(\alpha, \beta)$

The forward and backward variables  $(\alpha, \beta)$   $\alpha_{\tau}^{n,k}(i) = p(z_{n,\tau}^k = i | a_{n,0:\tau}^k, o_{n,1:\tau}^k, \Theta_n)$  and  $\beta_{t,\tau}^{n,k}(i) = \frac{p(a_{n,\tau+1:t}^k | z_{n,\tau}^k = i, o_{n,\tau+1:t}^k, \Theta_n)}{\prod_{\tau'=t}^{\tau} p(a_{n,\tau'}^k | h_{n,\tau'}^k, \Theta_n)}$ ,  $\forall n, k, t, \tau$  defined in section 3.2 are similar to the forward-backward messages in the Baum-Welch algorithm for hidden Markov models [31]. These variables are computed recursively by each agent using (34)-(36).

$$\alpha_{\tau}^{n,k}(i) = \begin{cases} \frac{\mu(z_{n,0}^k=i)\pi(z_{n,0}^k=i, a_{n,0}^k)}{p(a_{n,0}^k | h_{n,0}^k, \Theta_n)} & \text{if } \tau = 0 \\ \frac{\sum_{j=1}^{|\mathcal{Z}_n|} \alpha(z_{n,\tau-1}^k = j) W(z_{n,\tau-1}^k = j, a_{n,\tau-1}^k, o_{n,\tau}^k, z_{n,\tau}^k = i) \pi(z_{n,\tau}^k = i, a_{n,\tau}^k)}{p(a_{n,\tau}^k | h_{n,\tau}^k, \Theta_n)} & \text{if } \tau > 0 \end{cases} \tag{34}$$

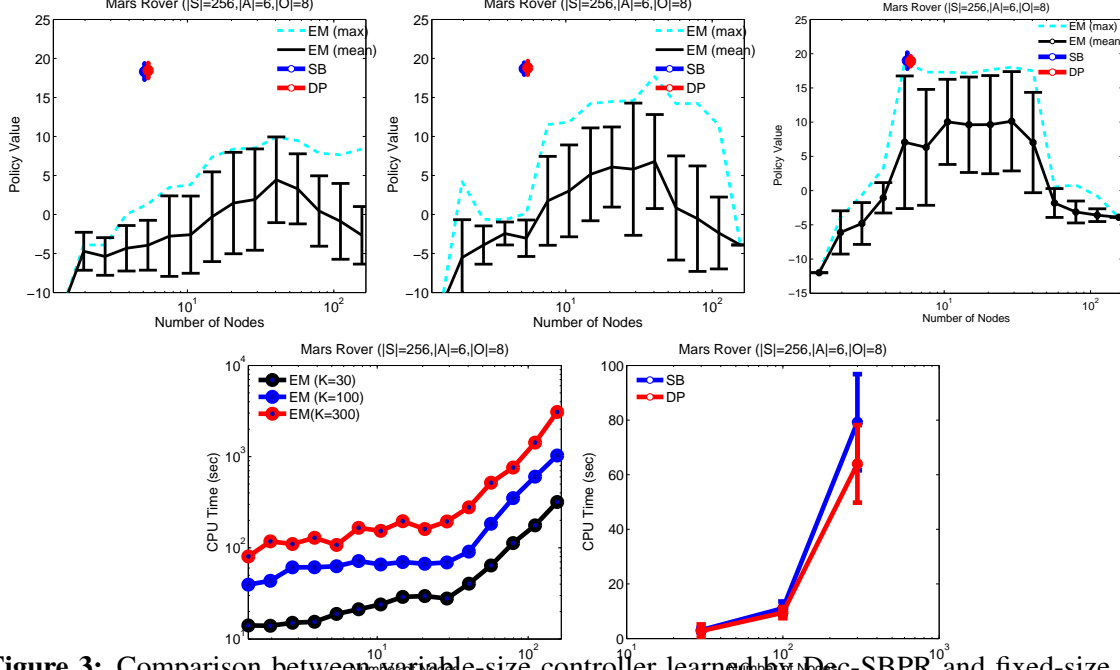
$$\beta_{t,\tau}^{n,k}(j) = \begin{cases} \frac{1}{p(a_{n,0}^k | h_{n,0}^k, \Theta_n)} & \text{if } \tau = 0 \\ \frac{\sum_{n=1}^{|\mathcal{Z}_n|} W(z_{n,\tau}^k = i, a_{n,\tau-1}^k, o_{n,\tau}^k, z_{n,\tau+1}^k = j) \pi(z_{n,\tau+1}^k = j, a_{n,\tau+1}^k) \beta_{n,t,\tau+1}^k(j)}{p(a_{n,\tau}^k | h_{n,\tau}^k, \Theta_n)} & \text{if } \tau > 0 \end{cases} \tag{35}$$

$$p(a_{n,\tau}^k | h_{n,\tau}^k, \Theta_n) = \begin{cases} \sum_{j=1}^{|\mathcal{Z}_n|} \mu(z_{n,0}^k = i) \pi(z_{n,0}^k = i, a_{n,0}^k) & \text{if } \tau = 0 \\ \sum_{i,j=1}^{|\mathcal{Z}_n|} \alpha(z_{n,\tau-1}^k = i) W(z_{n,\tau-1}^k = i, a_{n,\tau-1}^k, o_{n,\tau}^k, z_{n,\tau}^k = j) \pi(z_{n,\tau}^k = j, a_{n,\tau}^k) & \text{if } \tau > 0 \end{cases} \tag{36}$$

## D Additional Experimental results

### D.1 learning variable-size FSC vs learning fixed-size FSC

Additional experiments are added to study the impact of number of training samples.



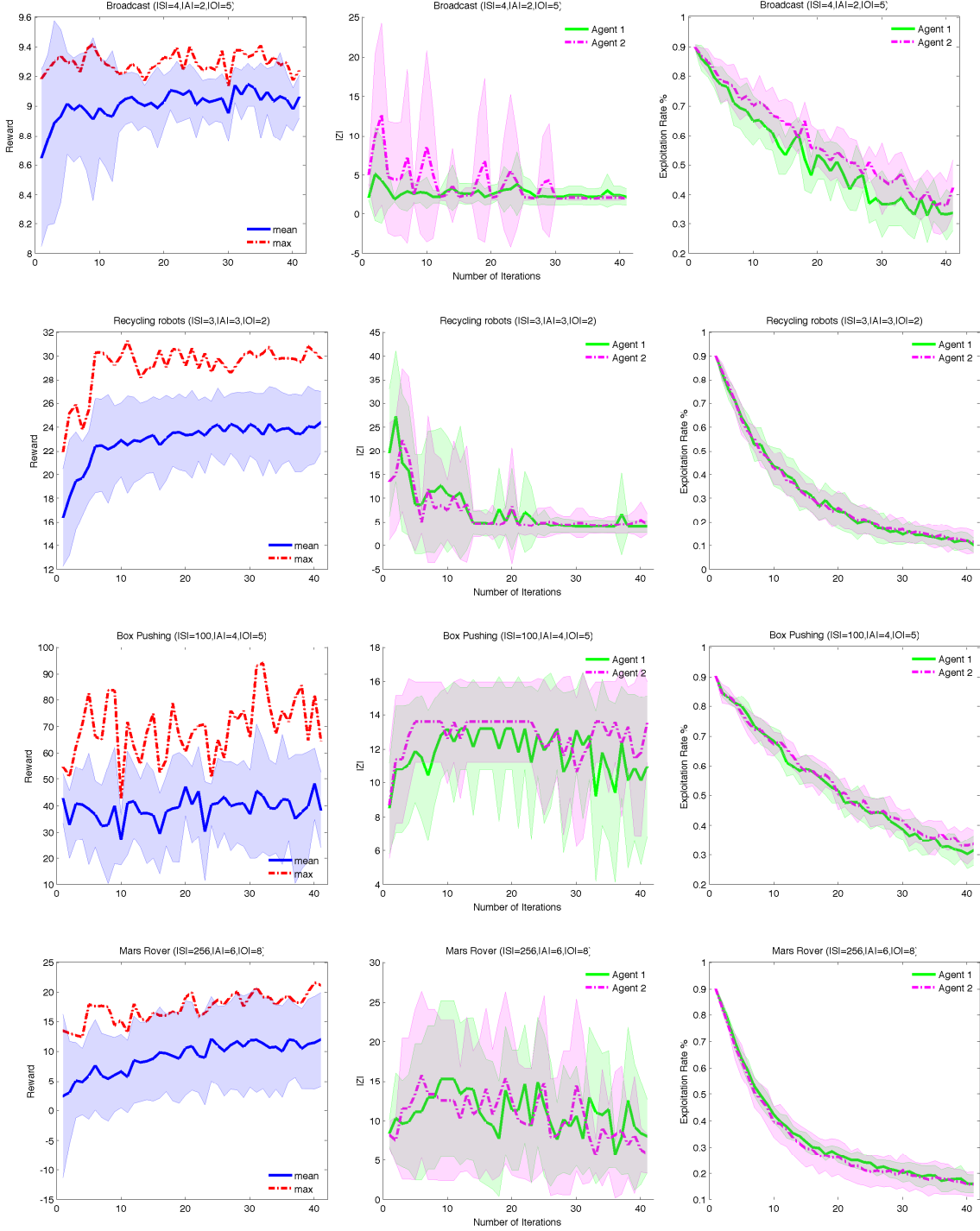
**Figure 3:** Comparison between variable-size controller learned by Dec-SBPR and fixed-size controller learned by EM algorithm. Top: testing value of policies learned using different number of training episodes ( $K = 30$ (left),  $K = 100$ (middle),  $K = 300$ (right)); Bottom: averaged computation time of EM (left) and SB (right).

## D.2 Sequential Batch Learning with Exploration Exploitation Trade-offs

Additional results from sequential batch learning with exploration exploitation trade-offs for five domains are plotted in Figure 4. In each iteration, a batch of samples are collected with updated behavior policies<sup>8</sup> and are used to learn a set new policies with Algorithm 1. These results are associated with the numbers reported in the first two columns of Table 2 in the main body of the paper.



<sup>8</sup>To generate these plots, 50 trajectories are collected in each iteration and the exploration parameter is set to be  $u = 100$ .



**Figure 4:** Additional plots for illustrating exploration-exploitation tradeoff, including testing value (left), inferred controller numbers (middle) and exploration rate (right). In each iteration, a batch of samples are collected with updated behavior policies and are used to learn a set new policies with Algorithm 1.

## References

- [1] Benchmarks. <http://rbr.cs.umass.edu/camato/decpomdp/download.html>. 2015.
- [2] C. Amato, D. S. Bernstein, and S. Zilberstein. Optimizing fixed-size stochastic controllers for POMDPs and decentralized POMDPs. *J. Autonomous Agents and Multi-Agent Systems*, 2(3):293–320, 2010.
- [3] C. Amato, G. Chowdhary, A. Geramifard, N. K. Ure, and M. J. Kochenderfer. Decentralized control of partially observable Markov decision processes. In *Proc. of the Conf. on Decision and Control*, 2013.
- [4] C. Amato, G. D. Konidaris, G. Cruz, C. A. Maynor, J. P. How, and L. P. Kaelbling. Planning for decentralized control of multiple robots under uncertainty. In *ICRA*, 2015.
- [5] C. Amato and S. Zilberstein. Achieving goals in decentralized POMDPs. In *Proc. of the 8th Int’l Conf. on Autonomous Agents and Multiagent Systems*, 2009.
- [6] B. Banerjee, J. Lyle, L. Kraemer, and R. Yellamraju. Sample bounded distributed reinforcement learning for decentralized POMDPs. In *Proc. of the 26th AAAI Conf. on Artificial Intelligence*, 2012.
- [7] M. J. Beal. *Variational algorithms for approximate Bayesian inference*. PhD thesis, University of London, 2003.
- [8] D. Bernstein, S. Zilberstein, R. Washington, and J. Bresina. Planetary rover control as a Markov decision process. In *6th Int’l Symposium on Artificial Intelligence, Robotics, and Automation in Space*, 2001.
- [9] D. S. Bernstein, C. Amato, E. A. Hansen, and S. Zilberstein. Policy iteration for decentralized control of Markov decision processes. *J. Artificial Intelligence Research*, 34:89–132, 2009.
- [10] D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein. The complexity of decentralized control of Markov decision processes. *Mathematics of Operations Research*, 27(4):819–840, 2002.
- [11] C. M. Bishop et al. *Pattern recognition and machine learning*, volume 1. springer New York, 2006.
- [12] M. Bryant and E. B. Sudderth. Truly nonparametric online variational inference for hierarchical dirichlet processes. In *Advances in Neural Information Processing Systems*, 2012.
- [13] C. Cai, X. Liao, and L. Carin. Learning to explore and exploit in POMDPs. In *Proc. of the 23th Annual Conf. on Neural Information Processing Systems*, 2009.
- [14] A. Dempster, N. Laird, and D. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *J. Royal Statistical Society B*, 39:1–38, 1977.
- [15] J. S. Dibangoye, O. Buffet, and F. Charpillet. Error-bounded approximations for infinite-horizon discounted decentralized POMDPs. In *Machine Learning and Knowledge Discovery in Databases*. Springer, 2014.
- [16] F. Doshi-Velez, D. Wingate, N. Roy, and J. B. Tenenbaum. Nonparametric Bayesian policy priors for reinforcement learning. In *Advances in Neural Information Processing Systems*, 2010.
- [17] T. S. Ferguson. A Bayesian analysis of some nonparametric problems. *Annals of Statistics*, 1:209–230, 1973.
- [18] H. Ishwaran and L. F. James. Gibbs sampling methods for stick-breaking priors. *J. Am. Statist. Assoc.*, 96(453), 2001.
- [19] A. Kumar and S. Zilberstein. Anytime planning for decentralized POMDPs using expectation. In *Proc. of the 26th Conf. on Uncertainty in Artificial Intelligence*, 2010.
- [20] H. Li, X. Liao, and L. Carin. Multi-task reinforcement learning in partially observable stochastic environment. *J. Machine Learning Research*, 10(1):1131–1186, 2009.

- [21] M. Liu, X. Liao, and L. C. and. Infinite regionalized policy representation. In *Proc. of the 28th Int'l Conf. on Machine Learning*, 2011.
- [22] M. Liu, X. Liao, and L. Carin. Online expectation maximization for reinforcement learning in POMDPs. In *Proc. of 23rd Int Joint Conf. on Artificial Intelligence*, 2013.
- [23] M. F. Martins and Y. Demiris. Learning multirobot joint action plans from simultaneous task execution demonstrations. In *Proc. of the 9th Int'l Conf. on Autonomous Agents and Multiagent Systems*, 2010.
- [24] J. Messias, M. Spaan, and P. Lima. Multi-robot planning under uncertainty with communication: a case study. In *Fifth Workshop on Multi-agent Sequential Decision Making in Uncertain Domains (MSDM)*, 2010.
- [25] F. A. Oliehoek. Decentralized POMDPs. In *Reinforcement Learning: State of the Art, Adaptation, Learning, and Optimization*, pages 471–503. Springer Berlin Heidelberg, Berlin, Germany, 2012.
- [26] J. Paisley and L. Carin. Hidden Markov models with stick-breaking priors. *Signal Processing, IEEE Trans. on*, 57(10):3905–3917, 2009.
- [27] J. Pajarinen and J. Peltonen. Periodic finite state controllers for efficient POMDP and DEC-POMDP planning. In *Proc. of the 25th Annual Conf. on Neural Information Processing Systems*, 2011.
- [28] O. Papaspiliopoulos and G. O. Roberts. Retrospective Markov chain Monte Carlo methods for Dirichlet process hierarchical models. *Biometrika*, 95(1):169–186, 2008.
- [29] J. Pineau, G. Gordon, and S. Thrun. Point-based value iteration: an anytime algorithm for POMDPs. In *Proc. of the 18th Int Joint Conf. on Artificial Intelligence*, 2003.
- [30] F. Qi and L. Losonczi. Bounds for the ratio of two gamma functions. *Journal of Inequalities and Applications*, 2010:204, 2010.
- [31] L. R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
- [32] C. P. Robert and G. Casella. *Monte Carlo statistical methods*, volume 319. Citeseer, 2004.
- [33] C. Wang and D. Blei. Variational Inference in Nonconjugate Models. *Journal of Machine Learning Research*, 14:1005–1031, 2013.
- [34] T.-T. Wong. Generalized dirichlet distribution in bayesian analysis. *Applied Mathematics and Computation*, 97(2):165–181, 1998.
- [35] F. Wu, S. Zilberstein, and N. R. Jennings. Monte-carlo expectation maximization for decentralized POMDPs. In *Proc. 23rd Int'l. Joint Conf. on Artificial Intelligence*, 2013.